

New developments in parsing Mizar

Czesław Bylinski¹ and Jesse Alama² *

Center for Artificial Intelligence
New University of Lisbon
Portugal
j.alama@fct.unl.pt

Abstract. The Mizar language aims to capture mathematical vernacular by providing a rich language for mathematics. From the perspective of a user, the richness of the language is welcome because it makes writing texts more “natural”. But for the developer, the richness leads to syntactic complexity, such as dealing with overloading.

Recently the Mizar team has been making a fresh approach to the problem of parsing the Mizar language. One aim is to make the language accessible to users and other developers. In this paper we describe these new parsing efforts and some applications thereof, such as large-scale text refactorings, pretty-printing, HTTP parsing services, and normalizations of Mizar texts.

1 Introduction

The Mizar system provides a language for declaratively expressing mathematical content and writing mathematical proofs. One of the principal aims of the Mizar project is to capture “the mathematical vernacular” by permitting authors to use linguistic constructions that mimic ordinary informal mathematical writing. The richness is welcome for authors of Mizar texts. However, a rich, flexible, expressive language is good for authors can lead to difficulties for developers and enthusiasts. Certain experiments with the Mizar language and its vast library of formalized mathematical knowledge (the Mizar Mathematical Library, or MML), naturally lead to rewriting Mizar texts in various ways. For some purposes one can work entirely on the semantic level of Mizar texts; one may not need to know precisely what the source text is, but only its semantic form. For such purposes, an XML presentation of Mizar texts has long been available [6]. However, for some tasks the purely semantic form of a Mizar text is not what is wanted. Until

* Supported by the ESF research project *Dialogical Foundations of Semantics* within the ESF Eurocores program *LogICCC* (funded by the Portuguese Science Foundation, FCT LogICCC/0001/2007). Research for this paper was partially done while a visiting fellow at the Isaac Newton Institute for the Mathematical Sciences in the program ‘Semantics & Syntax’. Karol Pąk deserves thanks for his patient assistance in developing customized Mizar text rewriting tools.

recently there has been no standalone tool, distributed with *Mizar*, that would simply parse *Mizar* texts and present the parse trees in a workable form.¹

Parsing texts for many proof assistants is often facilitated through the environment in which these proof assistants are executed. Thus, texts written for those systems working on top of a Lisp, such as *IMPS*, *PVS*, and *ACL2*, already come parsed, so one has more or less immediate access to the desired parse trees for terms, formulas, proofs, etc. Other systems, such as *Coq* and *HOL light*, use syntax extensions (e.g., *Camlp4* for Objective Caml) to “raise” the ambient programming language to the desired level of proof texts. For *Mizar*, there is no such ambient environment or read-eval-print loop; working with *Mizar* is more akin to writing a C program or *L^AT_EX* document, submitting it to *gcc* or *pdflatex*, and inspecting the results.

This paper describes new efforts by the *Mizar* team to make their language more workable and illustrates some of the fruits these efforts have already borne. This paper does not explain *how* to parse arbitrary *Mizar* texts. And for lack of space we cannot go into the detail about the *Mizar* system; see [3,4].

In Section 2, we discuss different views of *Mizar* texts that are now available. Section 3 describes some current applications made possible by opening up *Mizar* texts, and describes some HTTP-based services for those who wish to connect their own tools to *Mizar* services. Section 4 concludes by sketching further work and potential applications.

2 Layers of a *Mizar* text

It is common in parsing theory to distinguish various analyses or layers of a text, considered in the first place as a sequence of bytes or characters [1]. Traditionally the first task in parsing is **lexical analysis** or **scanning**: to compute, from a stream of characters, a stream of *tokens*, i.e., terminals of a production grammar *G*. From a stream of tokens one then carries out a **syntactic analysis**, which is the synthesis of tokens into groups that match the production rules of *G*.

One cannot, in general, lexically analyze *Mizar* texts without access to the MML. Overloading (using the same symbol for multiple, possibly unrelated meanings) already implies that parsing will be non-trivial, and overloading is used extensively in the *Mizar* library. Even with a lexical analysis of a *Mizar* text, how should it be understood syntactically? Through *Mizar*’s support for **dependent types**, the overloading problem is further complicated. Consider, for example, the *Mizar* fragment

```
let X be set,  
    R be Relation of X, Y;
```

¹ One parser tool, *lisppars*, is distributed with *Mizar*. *lisppars* is mainly used to facilitate authoring *Mizar* texts with Emacs [5]; it carries out fast lexical analysis only and does not output parse trees.

The notion of a (binary) relation is indicated by the non-dependent (zero-argument) type **Relation**. There is also the binary notion *relation whose domain is a subset of X and whose range is a subset of Y* , which is expressed as **Relation of X,Y** . Finally, we have the one-argument notion *relation whose domain is a subset of X and whose range is a subset of X* which is written **Relation of X** . In the text fragment above, we have to determine which possibility is correct, but this information would not be contained in a token stream (is Y the second argument of an instance of the binary **Relation** type, or is it the third variable introduced by the **let**?).

2.1 Normalizations of Mizar texts

One goal of opening up the Mizar parser is to help those interested in working with Mizar texts to not have to rely on the Mizar codebase to do their own experiments with Mizar texts. We now describe two normalizations of (arbitrary) Mizar texts, which we call weakly strict and more strict. The results of these two normalizations on a Mizar text can be easily parsed by a standard LR parser, such as those generated by the standard tool *bison*² and have further desirable syntactic and semantic properties. Other normalizations beyond these two are certainly possible. For example, whitespace, labels for definitions, theorems, lemmas, etc., are rewritten by the normalizations we discuss; one can imagine applications where such information ought not be tampered with.

2.2 Weakly strict Mizar

The aim of the weakly strict Mizar (WSM) transformation is to define a class of Mizar texts for which one could easily write an standard, standalone parser that does not require any further use of the Mizar tools. In a weakly strict Mizar text all notations are disambiguated and fully parenthesized, and all statements take up exactly one line. (This is a different transformation than single-line variant AUT-SL of the Automath system [2].) Consider:

```
reserve P,R for Relation of X,Y;
```

This Mizar fragment is ambiguous: it is possible that the variable Y is a third reserved variable (after the variables P and R), and it is possible that Y is an argument of the dependent type **Relation of X,Y** . The text becomes disambiguated by the weakly strict Mizar normalization to

```
reserve P , R for ( Relation of X , Y ) ;
```

and now the intended reading is syntactically evident, thanks to explicit bracketing and whitespace. (Any information that is implicitly contained by whitespace structure in the original text is destroyed.)

The result of the one-line approach of the weakly strict Mizar normalization is, in many cases, excessive parenthesization, unnecessary whitespace, and rather

² <http://www.gnu.org/software/bison/>

long lines.³ The point of the weakly strict Mizar normalization is not to produce attractive human-readable texts. Instead, the aim is to transform Mizar texts so that they have a simpler grammatical structure.

2.3 More Strict Mizar

A second normalization that we have implemented is called, for lack of a better term, more strict Mizar (MSM). The aim of the MSM normalization is to define a class of Mizar texts that are canonicalized in the following ways:

- From the name alone of an occurrence of a variable one can determine the category (reserved variable, free variable, bound variable, etc.) to which the occurrence belongs. (Such inferences are of course not valid for arbitrary Mizar texts.)
- All formulas are labeled, even those that were unlabeled in the original text.
- Some “syntactic sugar” is expanded.
- Toplevel logical linking is replaced by explicit reference. Thus,

```
 $\phi$ ; then  $\psi$ ;
```

using the keyword **then** includes the previous statement (ϕ) as the justification of ψ . Under the MSM transformation, such logical relationships are rewritten as

```
Label1:  $\phi$ ;  
Label2:  $\psi$  by Label1;
```

Now both formulas have new labels **Label1** and **Label2**. The logical link between ϕ and ψ , previously indicated by the keyword **then**, is replaced by an explicit reference to the new label (**Label1**) for ϕ .

- All labels of formulas and names of variables in a Mizar are serially ordered.

MSM Mizar texts are useful because they permit certain “semantic” inferences to be made simply by looking at the syntax. For example, since all formulas are labeled and any use of a formula must be done through its label, one can infer simply by looking at labels of formulas in a text whether a formula is used. By looking only at the name of a variable, one can determine whether it was introduced inside the current proof or was defined earlier.

3 Applications

Opening up the Mizar parser by providing new tools that produce parse trees naturally suggests further useful text transformations, such as pretty printing. An HTTP parsing service for these new developments is available for public

³ The longest line in the “WSM-ified” library has length 6042. About 60% (to be precise, 694) of the articles in the WSM form of the current version of the Mizar Mathematical Library (4.181.1147) have lines of length at least 500 characters. The average line length across the whole “WSM-ified” library is 54.7.

consumption. Four services are available. Submitting a suitable GET request to the service and supplying a Mizar text in the message body, one can obtain as a response the XML parse tree for the text, a pretty-printed form of it, or the WSM or MSM form of a text (either as plain text or as XML). The HTTP services permit users to parse Mizar texts without having access to the MML, or even the Mizar tools. See

<http://mizar.cs.ualberta.ca/parsing/>

to learn more about the parsing service, how to prepare suitable HTTP parsing requests, and how to interpret the results.

4 Conclusion and Future Work

Parsing is an essential task for any proof assistant. In the case of Mizar, parsing is a thorny issue because of the richness of its language and its accompanying library. New tools for parsing Mizar, with an eye toward those who wish to design their own Mizar applications without (entirely) relying on the Mizar tools, are now available. Various normalizations for Mizar texts have been defined. Further useful normalizations are possible. At present we are experimenting with a so-called “without reservations” Mizar (WRM), in which there are no so-called reserved variables; in WRM texts the semantics of any formula is completely determined by the block in which it appears, which should make processing of Mizar texts even more efficient.

References

1. Aho, A., Lam, M., Sethi, R., Ullman, J.: *Compilers: Principles, Techniques, and Tools*. Pearson/Addison Wesley (2007)
2. de Bruijn, N.G.: AUT-SL, a single-line version of Automath. In: Nederpelt, R., Geuvers, J.H., de Vrijer, R.C. (eds.) *Selected Papers on Automath, Studies in Logic and the Foundations of Mathematics*, vol. 133, chap. B.2, pp. 275–281. North-Holland (1994)
3. Grabowski, A., Kornilowicz, A., Naumowicz, A.: Mizar in a nutshell. *Journal of Formalized Reasoning* 3(2), 153–245 (2010)
4. Matuszewski, R., Rudnicki, P.: Mizar: the first 30 years. *Mechanized Mathematics and its Applications* 4(1), 3–24 (2005)
5. Urban, J.: MizarMode—an integrated proof assistance tool for the Mizar way of formalizing mathematics. *Journal of Applied Logic* 4(4), 414 – 427 (2006), <http://www.sciencedirect.com/science/article/pii/S1570868305000698>
6. Urban, J.: Xml-izing mizar: Making semantic processing and presentation of mml easy. In: Kohlhase, M. (ed.) *MKM. Lecture Notes in Computer Science*, vol. 3863, pp. 346–360. Springer (2005)